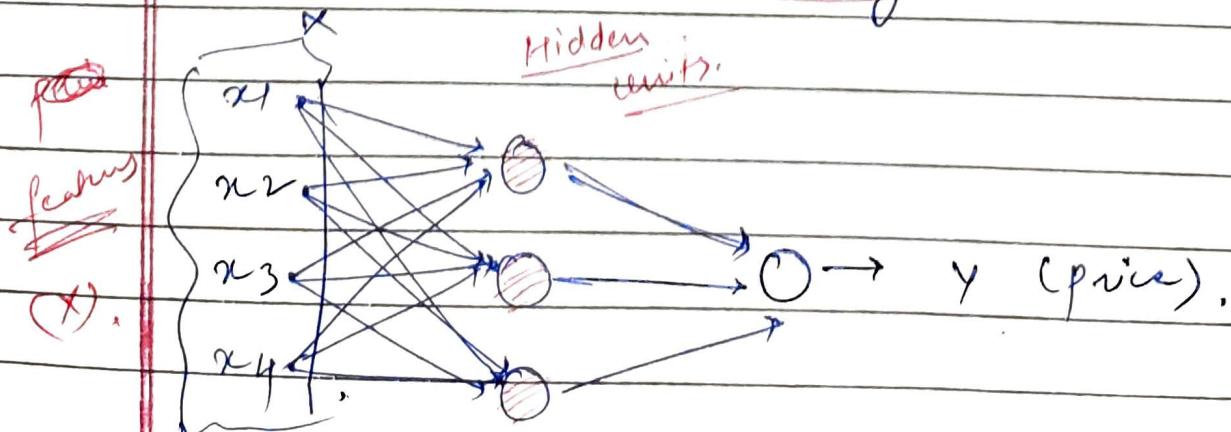
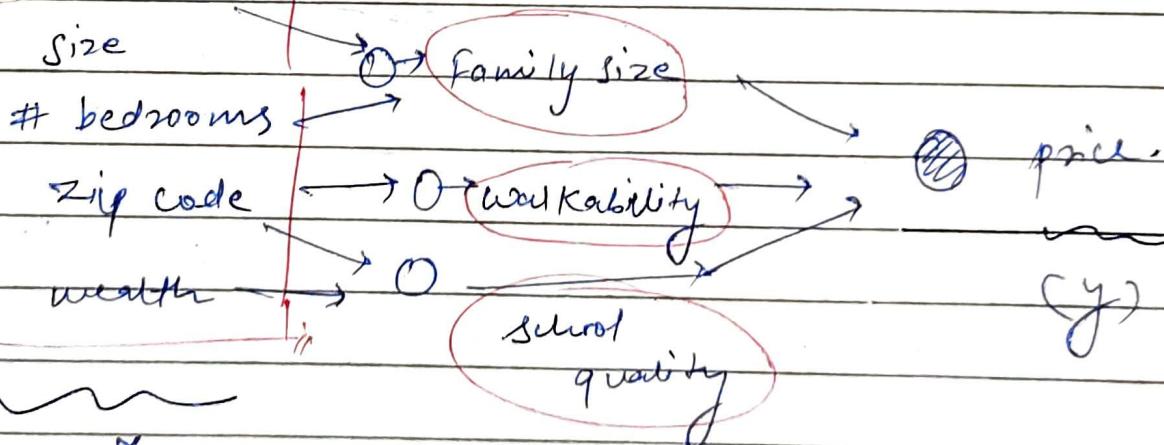
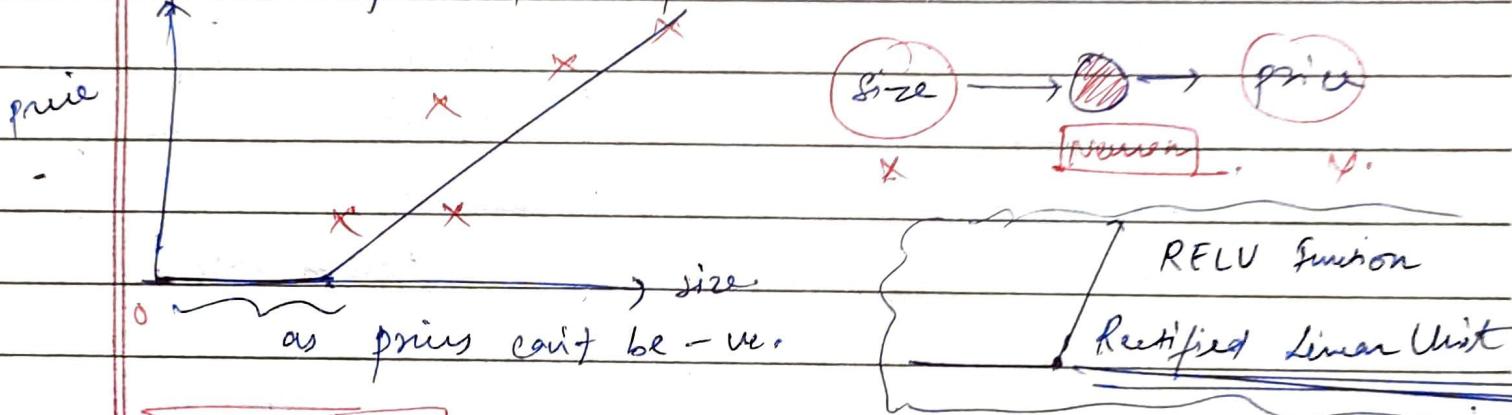


Course-1Neural network and deep learningby Andrew Ngweek 1INTRODUCTION TO DEEP LEARNINGWhat is Neural net

- deep learning refers to training neural netw. ex in housing price prediction.



Supervised Learning with neural networks

SUPERVISED LEARNING

(X)	(Y)	Application.
Home features	Price	Real Estate
Ad. clicking.	0/1	online Ad.
Image	object (1..10)	photo taggig
Audio	Text trans.	Speech recog
English	Chinese.	Machine transl.
Image / Radar info	position.	Autonomous driving. (Custom NN.)

why is deep learning taking off?

(d)

size of NN.

(e).

Traditional ml algorithm.

0

(m) Amount of data →

gradient ~ 0
learning being slow.

↑

ReLU

Deep Learning Specialization

- o ① Neural nw and deep learning.
- ② Improving deep neural nw.
- ③ Structuring your ML project.
- ④ Convolution neural network.
- ⑤ NLP: Building sequence model.

Page:

Date: / /

Week-2

NEURAL N/W BASICS

Binary Classification

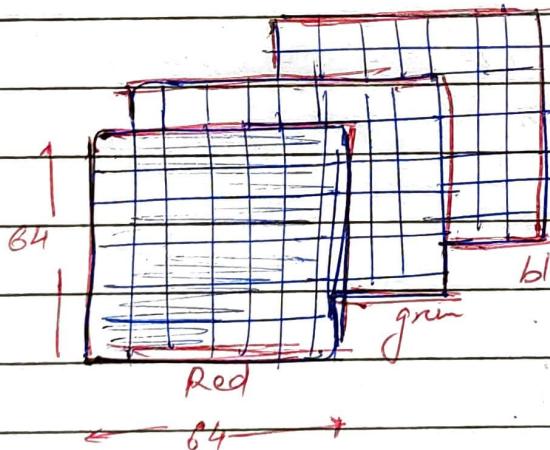


either x or $\neg x$

(1)

(0)

Storing an image.



To turn these pixel intensity values into feature vector, we row all feature intensity values in feature vector.

$$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \end{bmatrix} \text{ feature vector.}$$

* If image is 64×64

dimension of matrix $= (64 \times 64) \times 3 = 12288$.

$$n/n_x = 12288.$$

[# Notation]

$\rightarrow (x, y)$

$x \in \mathbb{R}^{n \times 1}, y \in \{0, 1\}$

$\rightarrow m$ training example.

$\rightarrow m_{\text{Train}} \rightarrow m_{\text{Test}}$

$$\mathbf{X} = \begin{bmatrix} & & & & \\ & x_1 & x_2 & \dots & x_m \\ & & & & \end{bmatrix} \quad \begin{matrix} \\ \\ \downarrow m \end{matrix}$$

$\leftarrow m \text{ columns} \rightarrow$

$$\boxed{\mathbf{X} \in \mathbb{R}^{n_x \times m}}$$

$$\mathbf{Y} = [y_1 \ y_2 \ \dots \ y_m]$$

$$\boxed{\mathbf{Y} \in \mathbb{R}^{1 \times m}}$$

Logistic Regression

Given \mathbf{x} , algo: $\boxed{\hat{y} = \text{prob.}(y=1 | \mathbf{x})}.$

$$\mathbf{x} \in \mathbb{R}^{n_x}$$

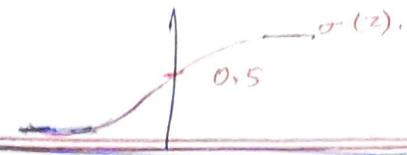
Parameters = $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$.

Output ?

$$\boxed{\hat{y} = w^T \mathbf{x} + b. \quad (\mathbf{x})}$$

$$\boxed{0 \leq \hat{y} \leq 1}$$

$$O/P = \hat{y} = \sigma(w^T x + b)$$



Page:

Date: / /

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$\left. \begin{array}{l} \text{if } z \text{ is very large } e^{-z} \approx 0 \\ \sigma(z) = \frac{1}{1+0} = 1 \end{array} \right\}$$

if z is << small, then $\sigma(z) = 1/(1+e^{-z}) \approx 0$.

4.18

Logistic Regression Cost function

$$\hat{y} = \sigma(w^T x + b), \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

given $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$.

$$z^{(j)} = w^T x^{(j)} + b$$

Loss function

$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

Want $\log \hat{y}$
large

$$y=1 = -\log \hat{y}$$

$$y=0 = \log(1-\hat{y})$$

\hat{y} = algorithm o/p

y = original label.

④ Loss function was defined for a single training example..

⑤ Cost function which measures over entire training set.

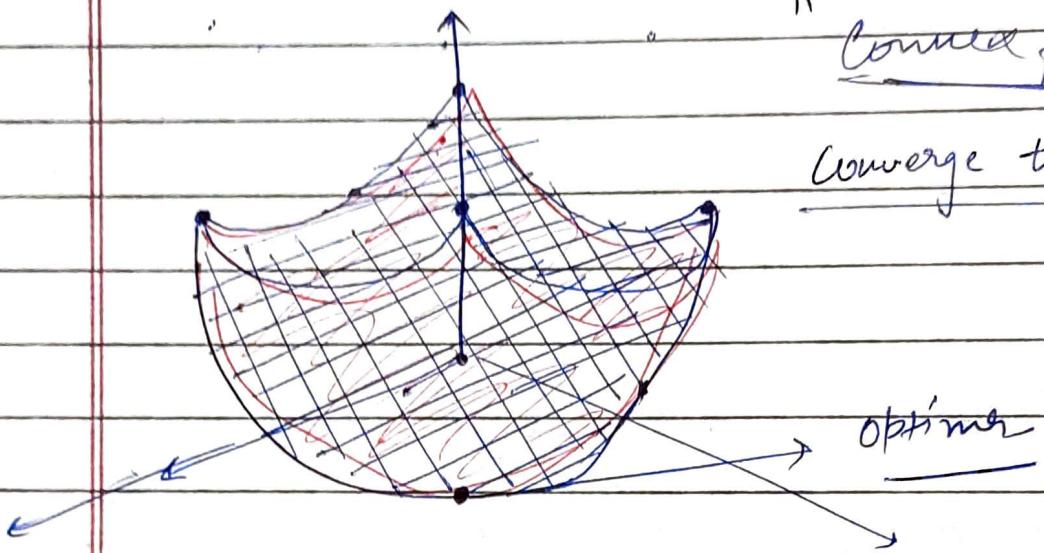
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$$

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m (y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_{\text{hat}(i)})).$$

Page: / /

Date: / /

Gradient Descent

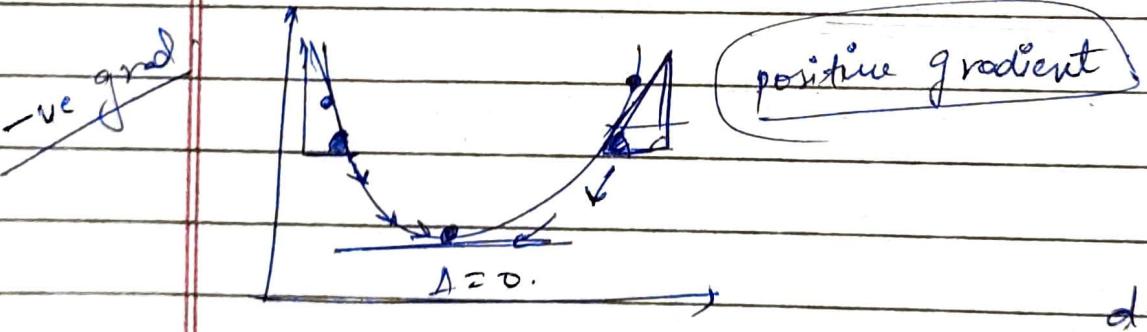


Converge to minima

Repeat until Convergence:

$$w := w - \alpha \frac{d J(w)}{d w}$$

learning rate.



$$\frac{\partial J(w, b)}{\partial w}$$

partial derivative

$$J(w, b) = \# w := w - \alpha \frac{\partial J(w, b)}{\partial w},$$

$$\# b := b - \alpha \frac{\partial J(w, b)}{\partial b},$$

∂b

Logistic Regression - Gradient descent

$$z = w^T x + b$$

$$\hat{y} = \sigma(z)$$

Page:

Date: / /

$$L(a, y) = -(y \log(a) + (1-y) \log(1-a)).$$

x_1

w_1

w_2

x_2

b

$$\frac{dL}{da} = \frac{d}{da} L(a, y) = -\frac{y}{a} + \frac{1-y}{1-a}.$$

$$\frac{dz}{dL} = \frac{dL}{dL} = \frac{\partial \log(\sigma)}{\partial a} \cdot \left(\frac{da}{dL} \right) = \frac{1}{a} - 1.$$

$$\frac{dL}{dw_1} = x_1 dz, \quad dL/dw_2 = x_2 dz, \quad db = dz.$$

gradient descent with m examples

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a_{(i)}, y_{(i)}).$$

$$a_{(i)} = \hat{y}_{(i)} = \sigma(z_{(i)}) = \sigma(w^T x_{(i)} + b).$$

$$J \geq 0; \quad dw = 0; \quad db = 0$$

for $i = 1 \rightarrow m \Rightarrow$

$$z_{(i)} = w^T x + b$$

$$a_{(i)} = \sigma(z_{(i)})$$

$$J+ = y_{(i)} \log(a_{(i)}) + (1-y_{(i)}) \log(1-a_{(i)})$$

$dL/dw_1 \text{ (i)}$

$dL/dw_2 \text{ (i)}$

$db \text{ (i)}$

$$dz_{(i)} = a_{(i)} - y_{(i)}$$

$$dw_1 + = x_{(i)} dz_{(i)}$$

$m = 8$

$$dw_2 + = x_{(i)} dz_{(i)}$$

$$db + = dz_{(i)}$$

$$d\omega_1 = \frac{dJ}{d\omega_1}$$

$$\omega_1 := \omega_1 - \alpha d\omega_1$$

Page:

Date: / /

$$\omega_2 := \omega_2 - \alpha d\omega_2$$

$$b := b - \alpha db.$$

Avoid explicit for loops and use Vectorization

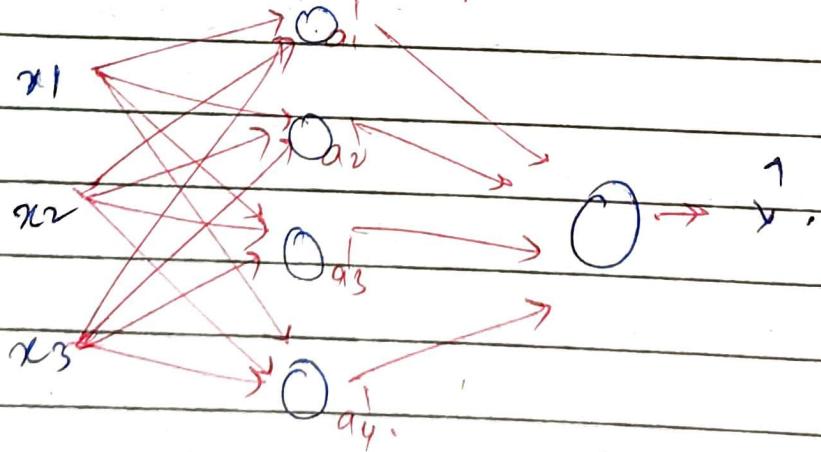
Vectorization

import numpy as np

$$z = np.dot(w, x) + b.$$

~~x — x — x — x — x —~~
week-3 Shallow Neural of w.

~~one hidden training layer~~



$$a = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix}$$

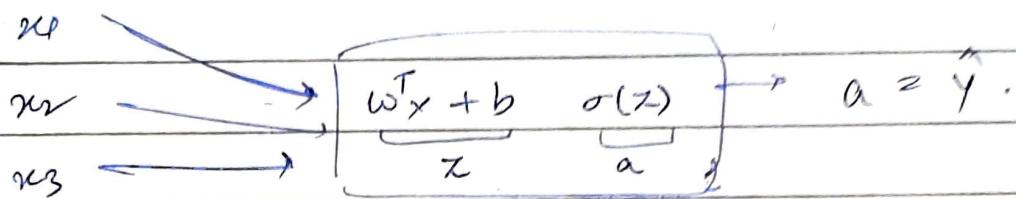
$0 \rightarrow$ hidden layer $\rightarrow 0 \rightarrow \hat{y}$.

Input
Layer

Page:

Date: / /

Computing a logistic regression NN Output \rightarrow



$$z^{(1)} = w_1^{(1)T} x + b_1^{(1)}$$

$$a_1^{(1)} = \sigma(z_1^{(1)})$$

$a_i^{(l)}$ \rightarrow target
 $a_i^{(l)}$ \rightarrow node.

Vectorized form

$$\begin{bmatrix} w_1^{(1)T} \\ w_2^{(1)T} \\ \vdots \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ \vdots \end{bmatrix}.$$

Vectorization over multiple examples.

$$z^{(1)(1)} = w^1 x^1 + b^1, \quad z^{(1)(2)} = w^1 x^2 + b^1,$$

$$z^{(1)(3)} = w^1 x^3 + b^1.$$

$$w^{(1)} = \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix} \quad w^1 x^1 = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad w^1 x^2 = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad w^1 x^3 = \begin{bmatrix} : \\ : \\ : \end{bmatrix}$$

$$w^{(1)} * \begin{bmatrix} x_1 & x_2 & x_3 & \dots \end{bmatrix} = \begin{bmatrix} z^{11} & z^{12} & z^{13} \end{bmatrix}$$

$$x = \alpha^{[0]}$$

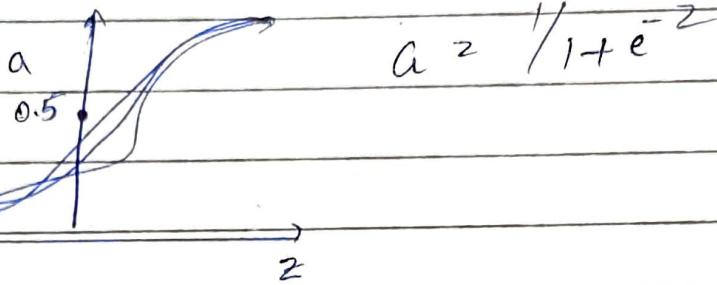
$$\underline{x}^{[1]} = w^{[1]} \alpha^{[0]} + b^{[1]}$$

Page:

Date: / /

Activation functions.

Sigmoid function



Tanh function

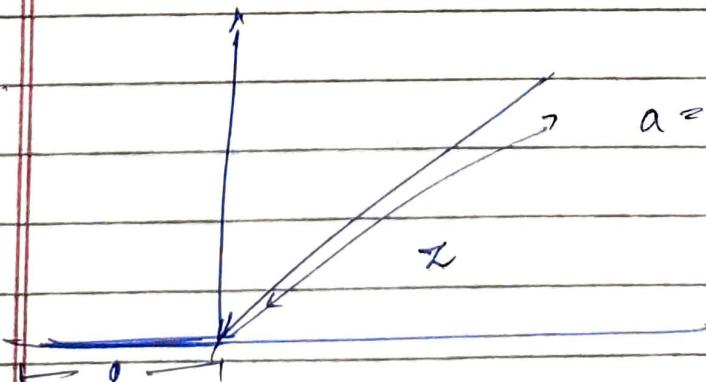
$$a = \tanh(z)$$

$$a = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \tanh(z).$$

* tanh function almost always superior than sigmoid function.

ReLU function

$$a = \max(0, z).$$



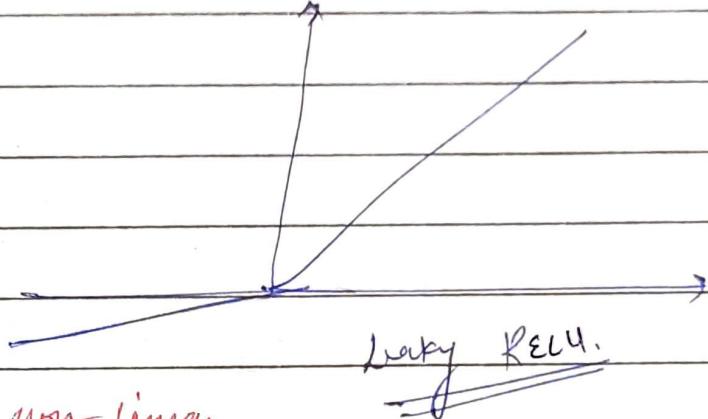
$$\tanh : a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Page:

Date: / /

$$\tanh(z) = \frac{\sinh z}{\cosh z} = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$z. \quad \frac{e^{2z} - 1}{e^{2z} + 1}$$



Why do we need non-linear activation function?

$$a^{(1)} = z^{(1)} = w^{(1)}x + b^{(1)}$$

$$a^{(2)} = z^{(2)} = w^{(2)}a^{(1)} + b^{(2)}$$

$$= w^{(2)} \cdot (w^{(1)}x + b^{(1)}) + b^{(2)}$$

$$= (w^{(2)} \cdot w^{(1)})x + (w^{(2)}b^{(1)} + b^{(2)})$$

$$= w'x + b'$$

$$g(z) \text{ or } \tanh(z) = 1 - \tanh(z)$$

Relu function = $(0, z < 0, ; 1, z \geq 0)$

$$\text{Sigmoid} = g(z)(1 - g(z))$$

gradient descent for neural net

Parameters = w^1, b^1, w^2, x^2

Cost function = $J(\text{Parameters})$

Page:

Date: / /

$$= \frac{1}{m} \sum_{i=1 \rightarrow n} L(y_i, \hat{y}_i). \quad \text{Backprop.}$$

gradient Descent \Rightarrow

repeat {

$$d\mathbf{w}^{(l)} = \frac{\partial J}{\partial \mathbf{w}^{(l)}}$$

$$d\mathbf{b}^{(l)} = \frac{\partial J}{\partial \mathbf{b}^{(l)}}$$

$$\mathbf{w}^{(l)} := \mathbf{w}^{(l)} - \alpha d\mathbf{w}^{(l)}$$

$$\mathbf{b}^{(l)} := \mathbf{b}^{(l)} - \alpha d\mathbf{b}^{(l)}$$

b^2

$$d\mathbf{z}^{(2)} = \mathbf{A}^{(2)} - \mathbf{Y}$$

$$d\mathbf{w}^{(2)} = \frac{1}{m} d\mathbf{z}^{(2)} \mathbf{A}^{(2)T}$$

$$d\mathbf{b}^{(2)} = \frac{1}{m} \text{np.sum}($$

$$d\mathbf{z}^{(1)}, \text{axis}=1, \text{inplace=True})$$

$$d\mathbf{z}^{(1)} = \mathbf{w}^{(2)T} d\mathbf{z}^{(2)} * g^{(1)}(-\mathbf{z}^{(1)})$$

$$d\mathbf{w}^{(1)} = \frac{1}{m} d\mathbf{z}^{(1)} \mathbf{x}^T$$

$$d\mathbf{b}^{(1)} = \frac{1}{m} \text{np.sum}(d\mathbf{z}^{(1)}, \text{axis}=1, \text{keepdims=True}),$$

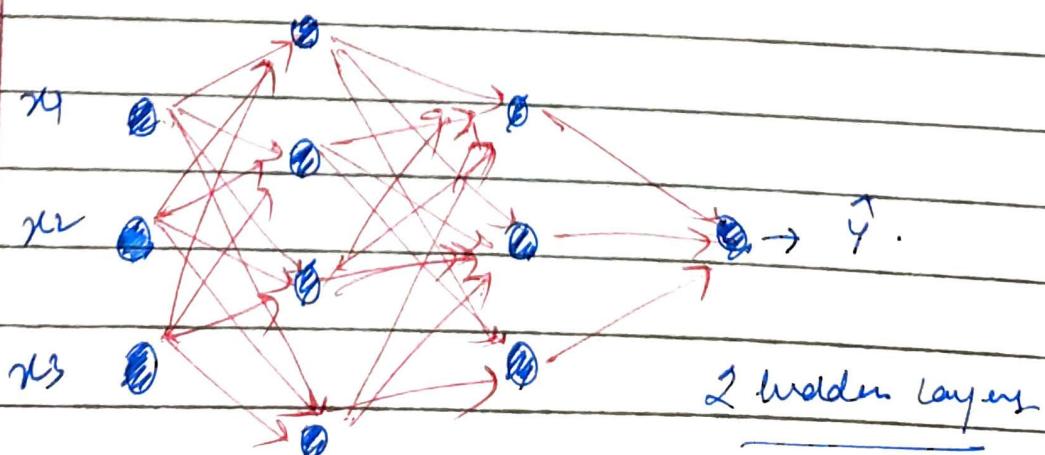
week-4

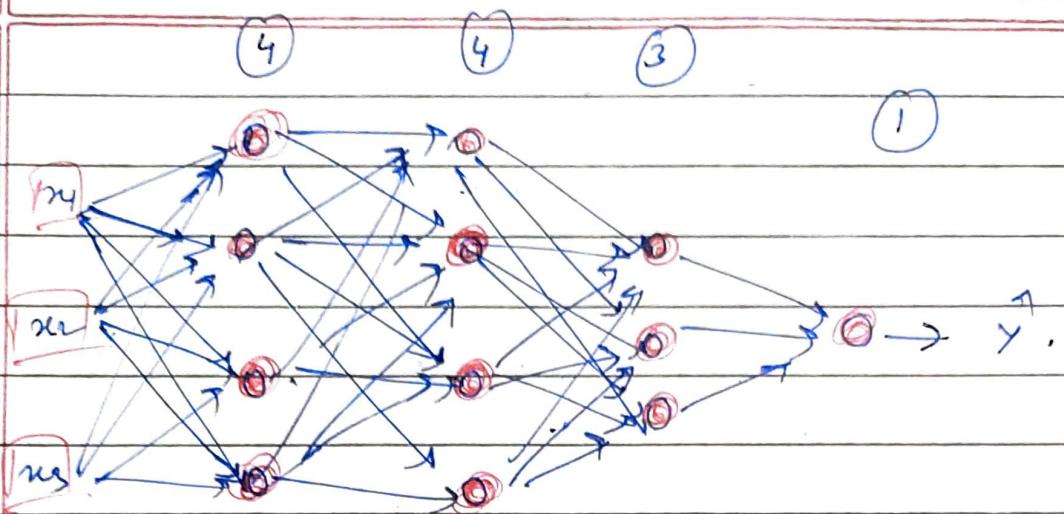
Deep Neural Networks.

(lecture)

Deep L Layer neural network

① Neural net with 2 hidden layers.





notations :-

l :- No. of layers.
 $a^{(l)}$:- activation in
 layer l .

$$a^{(l)} = g^{(l)}(z^{(l)})$$

$w^{(l)}$ = weights for $z^{(l)}$.
 $a^{(0)} = x$, $a^{(L)} = \hat{y}$.

$n^{(l)}$:- No. of units in
 layer l .

$$\begin{array}{c|c} n^{(1)} = 4 & n^{(4)} = 1 \\ n^{(2)} = 4 & n^{(0)} = n_x = 3 \\ n^{(3)} = 3 & \end{array}$$

Lesson 2 Forward propagation in Deep N/W

- For a single training example x .

$$x: \left[\begin{array}{l} z^{(1)} = w^{(1)} \cdot x + b^{(1)} \\ a^{(1)} = g^{(1)}(z^{(1)}) \end{array} \right] \quad l_1$$

$$\left[\begin{array}{l} z^{(2)} = w^{(2)} \cdot a^{(1)} + b^{(2)} \\ a^{(2)} = g^{(2)}(z^{(2)}) \end{array} \right] \quad l_2$$

$$z^{(4)} = w^{(4)} \cdot a^{(3)} + b^{(4)}$$

$$a^{(4)} = g(z^{(4)}).$$

$$= \hat{y}.$$

$$\boxed{x = a^0}.$$

$$\hat{A}^0 = x$$

(*)

$$z^{(L)} = w^{(L)} \cdot A^{(L-1)} + b^{(L)}$$

$$a^{(L)} = g(z^{(L)}).$$

Vectorized Version :-

for layer
n

l-... n
last time

explicit

for loop

$$z^{(L)} = w^{(L)} \cdot A^0 + b^{(L)}$$

$$A^{(L)} = g^{(L)}(z^{(L)}).$$

$$z^{(L)} = w^{(L)} \cdot A^{(L-1)} + b^{(L)}$$

$$A^{(L)} = g^{(L)}(z^{(L)}).$$

$$\hat{y} = g(z^{(L)}) = A^{(L)}$$

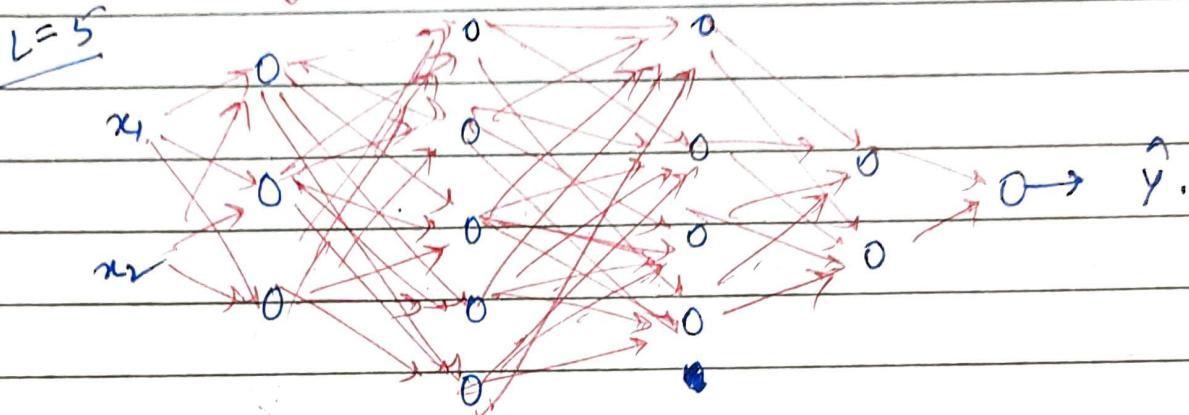
@ 6:15 $a^L = g(w, 2^W)$.
 $a, z \rightarrow (n^L, 1)$.

Page:

Date: / /

Review 3:

Getting your matrix dimensions right.



$$\bar{z}^{(1)} = w^{(1)}_i x + b^{(1)}$$

$(3,1) \quad (3,2) \quad (2,1)$
 $(n^{(1)}, 1), (n_1, n_0), (n_0, 1)$

$$\begin{array}{|c|c|} \hline n^{(1)} & n^{(3)} \\ \hline n^{(2)} & n^{(4)} \\ \hline n^{(5)} & n^0 \\ \hline \end{array} \begin{array}{l} = 3 \\ = 5 \\ = 1 \end{array} \begin{array}{l} = 4 \\ = 2 \\ = n_x = 2 \end{array}$$

$$\begin{aligned} \bar{z}^{(1)} &= (3,1) = (n^{(1)}, 1) = a^{(1)} \\ \rightarrow w^{(1)} &= (3,2) = (n_1, n_0) \Rightarrow (n_1, n_{(1)}) \\ x &= (2,1) = (n_0, 1) \end{aligned}$$

for dimension b :

$$\rightarrow b = (3,1) \Rightarrow (n^{(1)}, 1) \text{ and}$$

$$\rightarrow dw \approx w, db \approx b$$

w and \bar{z} have same dims.



$$\underline{z}^{[l]} = \underline{w}^{[l]} \cdot \underline{x} + \underline{b}^{[l]}$$

$(n, 1) \quad (n, n^0) \quad (n^0, 1) \quad (n^{[l]}, 1)$

$$\underline{z}^l = \begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \\ | & | & \dots & | \\ 1 & 1 & \dots & 1 \end{bmatrix} \quad (n, m)$$

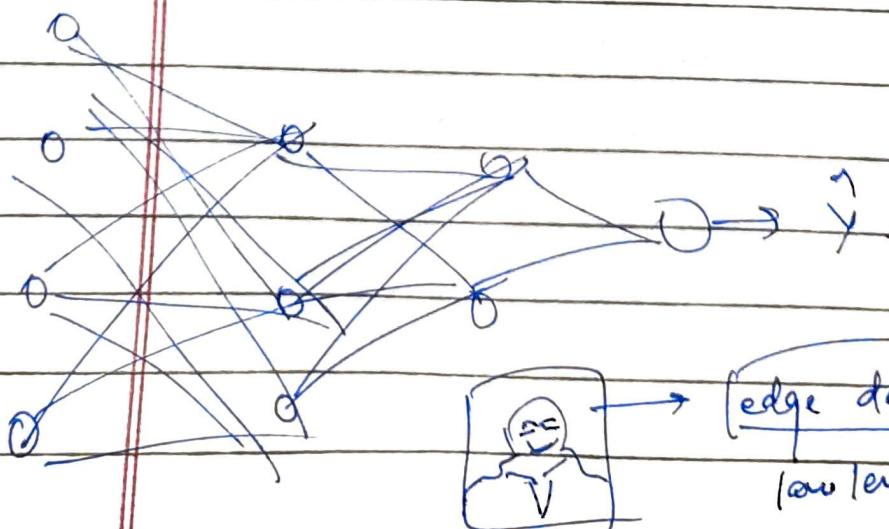
$$\underline{z}^{(l)}, a^{(l)} : (n^l, 1)$$

* given $\underline{z}^l, A^{(l)} : (n^l, m)$

if $l=0$ $A^{(0)} = X = (n_0, m)$

$$d\underline{z}^l, dA^{(l)} : (n^l, m)$$

Lecture 7: Why deep representations?



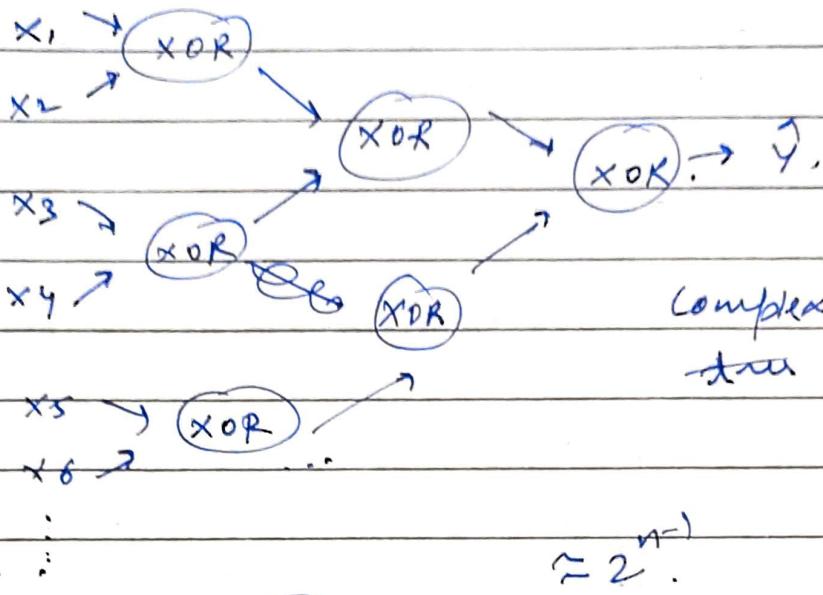
Circuit theory and deep learning

Ex: $y = x_1 \text{ XOR } x_2 \text{ XOR } x_3 \text{ XOR } x_4 \dots$

Page:

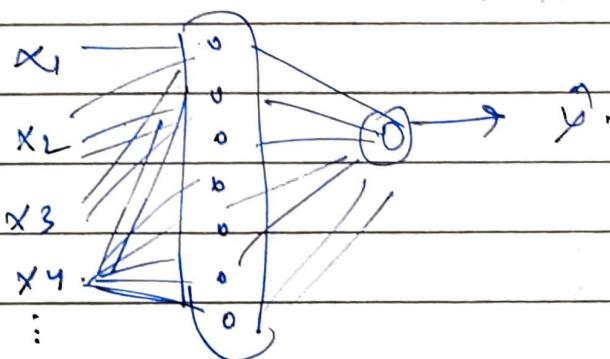
Date: / /

deep
neural
network



Complexity : $O(\log n)$.

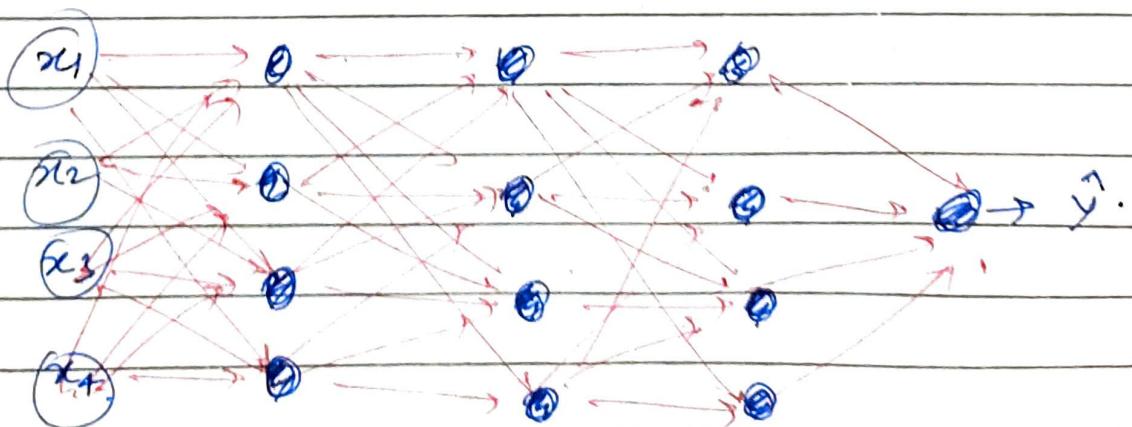
single
layer
 N/N :



Complexity : 2^n
exponentially \gg large

Lesson - 5

Building blocks of deep N/N



for Layer L : $w^{(L)}$, $b^{(L)}$

forward prop. : $a^{(L-1)}$, o/p : $a^{(L)}$

$$z^{(L)} = w^{(L)} \cdot a^{(L-1)} + b, \text{ calc } z^{(L)}$$

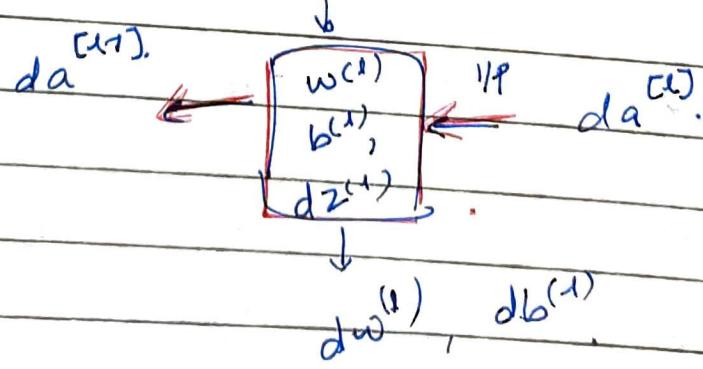
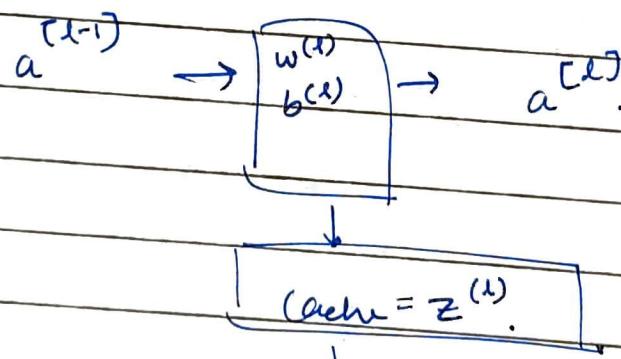
$$a^{(L)} = g(z^{(L)}).$$

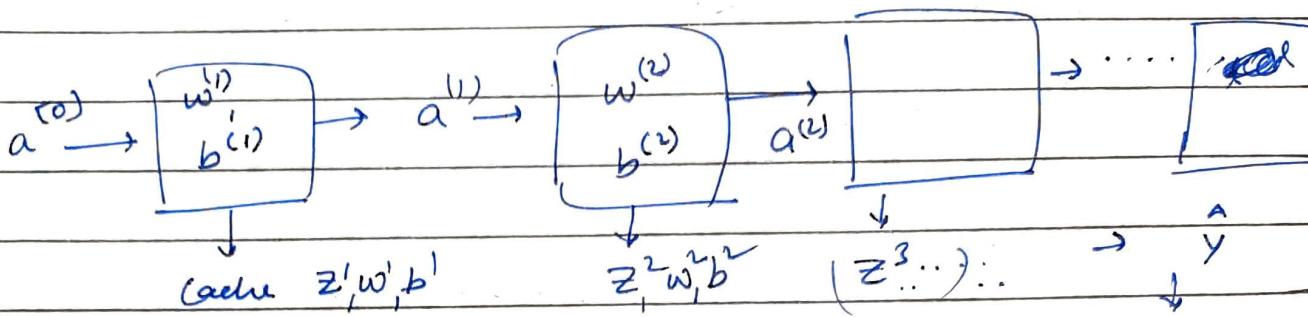
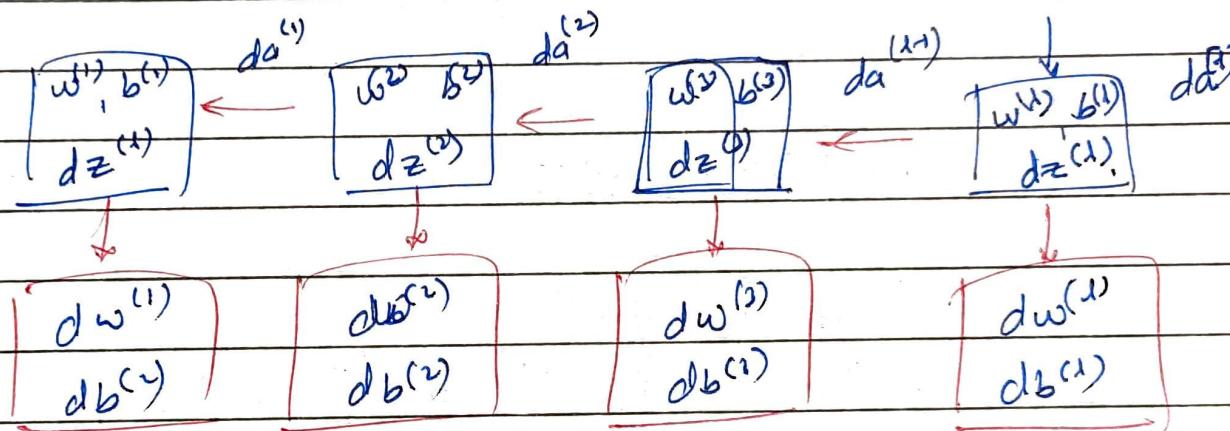
Backward prop. \Rightarrow

I/P : $da^{(L)}$ o/p : $da^{(L-1)}$.
 (calc (z^W)). dw^W

Summary

layer L.



forward prop.backward prop.

$$w^{(L)} := w^{(L)} - \alpha \cdot d\omega^{(L)}$$

$$b^{(L)} := b^{(L)} - \alpha \cdot db^{(L)}$$

Learn 6.~~Parameters~~ \neq Hyperparametersforward \neq back propagation.FORWARD PROP.

→ Input $a^{[L-1]}$

→ Output $a^{(L)}$, cache $z^{(L)}$

$$z^{(l)} = w^{(l)} \cdot a^{(l-1)} + b^{(l)}$$

$$a^{(l)} = g^{(l)}(z^{(l)}).$$

Vectorized:

$$\mathbf{z}^{(l)} = \mathbf{w}^{(l)} \cdot \mathbf{A}^{(l-1)} + \mathbf{b}^{(l)}$$

$$\mathbf{A}^{(l)} = g^{(l)}(\mathbf{z}^{(l)}).$$

BACKWARD PROP.

$$\rightarrow I/P : da^{(l)}$$

$$O/P : da^{(l-1)}, dw^{(l)}, db^{(l)}$$

$$dz^{(l)} = da^{(l)} * g^{(l)'}(z^{(l)}).$$

$$dw^{(l)} = dz^{(l)} \cdot a^{(l-1)}.$$

$$db^{(l)} = dz^{(l)}.$$

$$da^{(l-1)} = w^{(l)\top} \cdot dz^{(l)}.$$

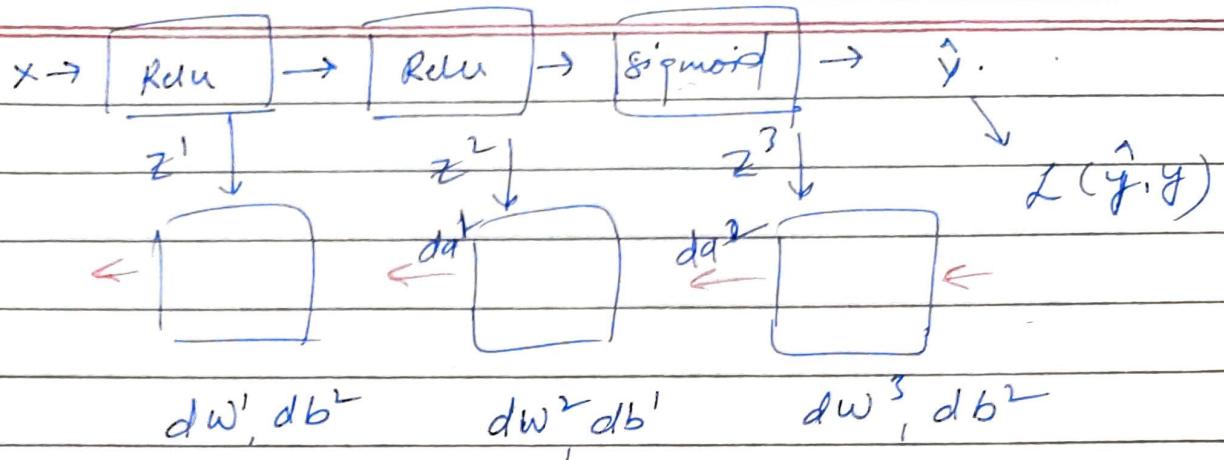
vector's

$$dz^{(l)} = da^{(l)} * g^{(l)'}(z^{(l)}).$$

$$dw^{(l)} = \frac{1}{m} dz^{(l)} \cdot A^{(l-1)\top}$$

$$db^{(l)} = 1/m \text{ np.sum } (dz^{(l)}, \text{axis} 2, \text{keepdim} z)$$

$$da^{(l-1)} = w^{(l)\top} \cdot dz^{(l)}.$$



$$da^{(L)} = \frac{-y}{a} + \frac{(1-y)}{(1-a)}$$

lecture 7

Parameters \neq Hyperparameters

PARAMETERS: $w^{(0)}, b^{(0)}, w^{(1)}, b^{(1)}, \dots, w^{(n)}, b^{(n)}$

HYPERPARAMETERS :- learning rate (α), #iterations,
#hidden layer / units, activation fx.

$dz^{(L)}$ $= A^{(L)} - y$.

$$dW^{(L)} = \frac{1}{m} dz^{(L)} \cdot A^{(L-1)}^T$$

$$db^{(L)} = \frac{1}{m} \text{np.sum}(dz^{(L)}, \text{axis}=1, \text{keepdim=True})$$

$$dz^{(L-1)} = W^{(L)} \cdot g'(z^{(L)})$$